# On Line Appendices

# Getting the Foot out of the Pelvis:
## Modelling Problems affecting Use of  SNOMED-CT Hierarchies in Practical Applications

**Alan Rector, Thomas Schneider, Sam Brandt**
**JAMIA vol 18, 2011**

# Table of Contents

## Notes

### Notation & diagramming conventions

For demonstrating and explaining the model, we have used a slightly abbreviated form of the Manchester OWL Syntax [1], in which we have replaced "Equivalent Classes" with "==" and SubClassOf with "→" and omitted the keywords CLASS and PROPERTY where they can be assumed. In the body of the text, we have used shortened forms of SNOMED preferred names. Where role groups have only a single member, they have usually been omitted for clarity in the body of the text, but the complete expression included in a footnote. (See also Appendix V: Overview of Description Logics).

In figures that are screen shots from Protégé 4, boldface indicates classes that have been introduced or changed in the repair; circles containing '≡' indicate fully defined concepts; open circles indicate partially defined concepts. In OWLViz graphs, orange ovals indicate fully defined concepts; yellow ovals represent partially defined concepts.

### Repetition of material in printed version

We have attempted to make each appendix comprehensible on its own without reference to the original printed version of the paper. Where appendices expanded on the material in the printed version, this inevitably involves some repetition of that material.

### Screen shots and figures in printed version

The printed version uses hand-drawn high-resolution diagrams. For documentation, this on-line supplement uses the corresponding screen shots from which the diagrams in the printed version were derived plus some additional screen shots documenting assertions made in the printed document.

## Appdendix I.      Details of modelling options for Myocardial Infarction

SNOMED does not classify "Myocardial infarction"[1] as a kind of "Ischemic heart disease"[2] although all references and experts consulted do. On investigation, the problem is that 'Myocardial infarction", "Infarction" and "Ischemic heart disease" are all only partially defined, so that there is no logical connection between them. Hence, the classifier cannot infer the linkage. The root problem is that "Infarction"[3] is not defined as being due to "Ischemia"[4], contrary to reference works and collaborators.

To solve this problem, all three concepts must be fully defined. However, this is not entirely straightforward.

In the original SNOMED formalism, there was a construct known as "right identities" that allowed chains of attributes to imply a single attribute. Most modern description logics including the SNOROCKET classifier, and OWL2 support a more general construct known as "property paths" that allows an arbitrary sequence of attributes to imply an attribute. Given this construct, we could state, for example, that the chain "morphology"[5] followed by "due to"[6] implies "due to". In OWL, the operator "o" is used between properties to indicate a chain, so this would appear as:

> morphology o due_to → due_to[7]

(This would extend the semantics of the existing attribute for "due to", so that it might be preferable to use a sub-property or even an entirely different property, but that detail can be ignored for now.)

In addition, we have to add a fully defined term roughly:

> "Disorder due to ischemia" == "Disorder" that "due to" some "Ischemia"[8]

In addition, we must modify the definition of "Ischemic heart disease" so that it is a fully defined term:.

> "Ischemic heart disease" == "Heart disease" and "Disorder due to ischemia"[9]

A second alternative may seem more direct to some users. Unlike the original SNOMED formalism, the SNOROCKET classifier also supports the combination of fully defined concepts with additional necessary conditions (subclass axioms) – a combination known in the description logic community as 'General Concept

---

[1] Myocardial infarction (disorder) | 22298006
[2] Ischemic stroke (disorder) | 422504002
[3] Infarct (morphologic abnormality) | 55641003
[4] Ischemia (disorder) | 52674009
[5] Associated morphology (attribute) | 116676008
[6] Due to (attribute) | 42752001
[7] 'Associated morphology' o 'Due to (attribute)' → 'Due to (attribute)'
[8] Ischemia (disorder) | 52674009
[9] 'Heart disease (disorder)' and "Disorder due to ischemia (disorder)"

Inclusion axioms' (GCIs). General Concept Inclusion axioms were initially thought to make reasoning intractable, but were found so useful in GALEN's GRAIL language [2] that they have since been extensively studied [3]. It has been shown that they can be included in the description logic fragment used by SNOMED (known as EL++ or OWL-EL) without seriously affecting performance.

In this case, General Inclusion Axioms provide an alternative construct with which SNOROCKET can deal. We first define ischemic disease as above, but then add a new class "Disorder due to Infarct"

> "Disorder due to Infarct" == Disease that "associated morphology" some Infarct[10]

Finally, we add the additional General Concept Inclusion Axiom

> "Disorder due to Infarct" → "Disorder due to Ischemia"

This construct gives the same hierarchy after classification, although the semantics of the two constructs are not strictly equivalent. (It is not true, in general, that General Concept Inclusion axioms can replace property paths or vice versa.)

In either case, if one makes "Myocardial ischemia" fully defined and asserts that "Angina" has the site myocardium rather than just heart, then one obtains the more radically rearranged but more compact structure in right hand half of Fig 2. Which is more appropriate is a matter for discussion amongst medical experts.
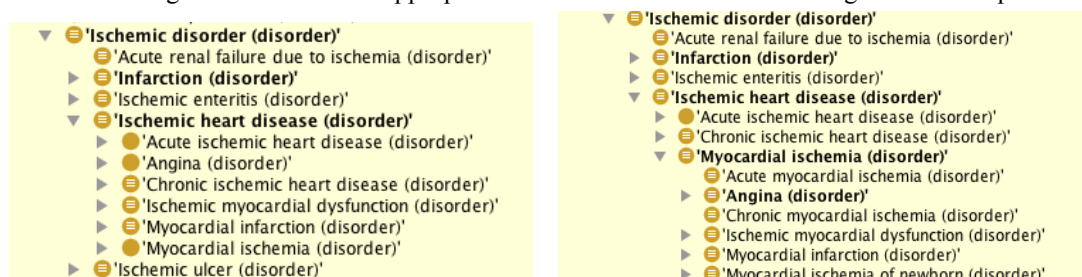


**Figure 2:Hierarchy under "Disorder due to Ischemia" after two different levels of changes (screen shots that form basis for Figure 2 in from printed paper for reference)**

## Appdendix II.        Modelling of branches so as to preserve transitivity

SNOMED currently models branches analogously to parts as children of the "Structure" (S) node in the SEP triple. This gives rise to serious errors in the classified hierarchies. One way to solve the problem is to move the branches to be children of the "Entire" (E) node in the SEP triple, but this sacrifices the transitivity of the branch_of relations, *i.e.* that branches of branches of a structure are also branches of that structure.

There are two solutions to this problem within EL++, both of which have been shown to work with the SNOROCKET classifier.

- We can declare a new attribute, "Is branch of (attribute)" and make it transitive. Then branches can be linked to the Entire (E) node using a new attribute instead of being made subclasses of the "Structure" (S) node. This is by far the simplest solution.

- Alternatively, the SEP triples can be extended to quintuples, with a new "Structure or Branch" (SB) as a parent of the current "Structure" (S) node and a new "Branch" (B) node as a child of the "Structure or Branch" (SB) node.

At the time SNOMED chose to adopt SEP triples, transitive properties resulted in significant decrease in performance of the then available classifiers. New algorithms developed since that time have largely eliminated this penalty, [4] so that there would seem to be little reason to adopt the second alternative with its proliferation of nodes. Indeed, it can be argued that classifier performance is no longer an argument for the SEP triple construct itself, which might therefore be dispensed with. In our experience, the distinction between the three nodes causes significant confusion to users.

## Appdendix III.        Broad and Narrow definitions of soft tissue

Some authorities define "soft tissues" as any tissue that is not bone. Other authorities are more specific and list only skin, fat, muscle, connective tissue, and organ tissues etc. but do not mention nerves, blood vessels or lymphatic vessels. Discussions with doctors and consulting text book tables of contents on what counts as disorders and injuries of soft tissues tends to support the second view, but not entirely consistently.

---

[10] 'Disease (disorder)' that RoleGroup some ('Associated morphology (attribute)' some 'Infarct (morphologic abnormality)')

SNOMED uses a relatively broad definition, including vessels and nerves but not internal organs. However, this broad definitions leads to a number of classifications that made experts uncomfortable, particularly that "Hypertension", "Neuropathy"[11]and "Arteritis"[12] etc. are classified as "Disorders of soft tissue"[13] Combined with the omission of "Skin and subcutaneous tissue" the unmodified list of "Disorders of soft tissues" appeared both over-long and insufficient. (See Figure Appe.III.1 ) Since there is a case for each concept, it seems useful to provide both concepts with distinct and unambiguous fully defined names – *e.g.:*

- "Soft (non-hard) tissues"
- "Soft tissues (not including vessels, nerves)"

Disorders and injuries referencing these concepts could then be named explicitly and correspondingly, and authors and users could select the concept that best fits their purpose.

The original hierarchy for "Disorder of soft tissue" is shown on the left in figure Appe.II.1; the alternative with both broader and narrower definitions of soft tissue and the addition of "Skin and subcutaneous tissue under "Soft tissue not including vessels and nerves) is shown on the right. The version on the right also includes "Skin and subcutaneous tissue under "Disorder of soft tissue (not including vessels)". As a result the individual concepts "Disorder of the skin of the head", "Disorder of the skin of the neck", etc. from the left hand side have been subsumed under "Disorder of the Skin AND/OR subcutaneous tissue".

For the purposes of our application, the version of "Disorder of soft tissue (not including vessels and nerves)" is more useful, although significant additional editing is needed under some of the subheadings. In addition, some collaborators pointed out that both versions of soft tissue exclude abdominal viscera, although thought it probably more of an issue of naming than of substance. The specific issue is whether or not rupture of the spleen should be considered as a soft tissue injury, under either broad or narrow definition. Most collaborators concurred with excluding viscera but thought that the fully defined name should make this clear, *e.g.* "Soft (not hard) tissue, not including viscera".
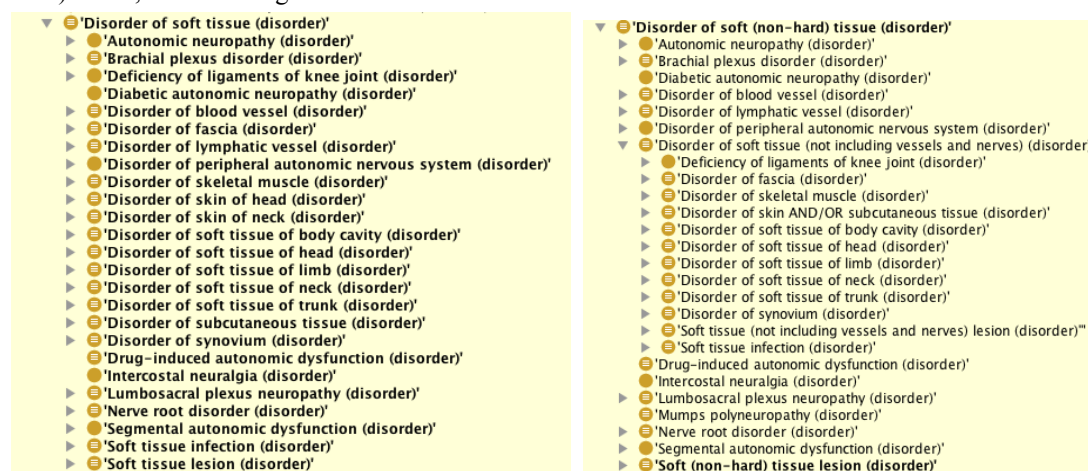


**Figure Appe.III.1: Disorders of soft tissue: original on left, modified version with inclusion of skin and subcutaneous tissues and addition of separate categories for "Soft (non-hard) tissue " and "Soft tissue (not including vessels and nerves)**

## Appdendix IV.    Issues in modelling the Endocrine System

The issue of the site of Diabetes mellitus is discussed in the main paper under the heading 'Is "Diabetes mellitus type 1" a "Disorder of the abdomen"?'. Similar issues arise with many endocrine disorders. Should "Thyrotoxicosis" be considered a "Disorder of the neck" or even necessarily of the thyroid gland? Should non-functional tumors of endocrine organs be considered endocrine disorders? What about dysregulation that involves several organs in the complex feedback loops involved in endocrine regulation?

There appears to be a case for a consistent pattern that distinguishes structural, functional, and regulatory disorders. There is also, clearly, a need for a parent concept to represent the disjunction and child concepts to represent the disjunction and conjunction of the three, possibly in various combinations.

---

[11] Neuropathy (disorder) | 386033004
[12] Arteritis (disorder) | 52089001
[13] Disorder of soft tissue (disorder) | 19660004

We suggest these issues are worth experiment and debate. It is difficult to decide a priori without the evidence of alternative attempts and their consequences via classification on the hierarchies. The full potential hierarchy is sketched in Figure Appe.IV.1. (Note that it is difficult to imagine a structural and regulatory disorder that was not also functional.) Not all of the concepts are relevant to all endocrine hormones or organs. There are several alternative modelling patterns that might be used to achieve similar results. Experiments are needed to show which gives the best results, whether such a uniform structure is practical, and if its benefits outweigh its costs.

An added complication is that the resulting structures would also have to be harmonised with SNOMED's distinction between observables and disorders, which is beyond the scope of this paper.
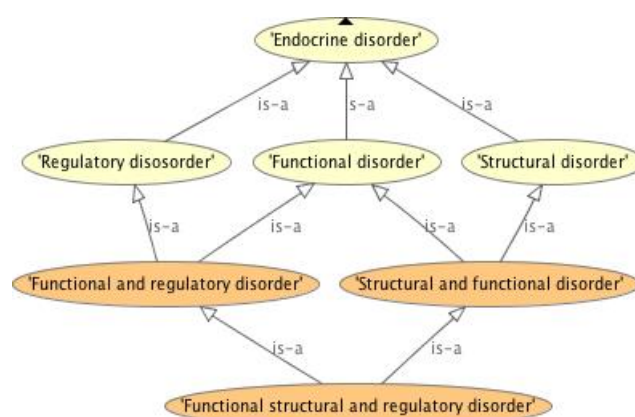


**Figure Appe IV.1: Sketch of potential hierarchy of structural, functional, and regulatory disorders in the endocrine system.**

# Appdendix V. Brief overview of description logic semantics

SNOMED-CT is formulated in a description logic. Description logics emerged in the 1990s as the dominant form of structured knowledge representation for terminologies and "ontologies". Recently, syntaxes and APIs for description logics have been standardised in the Web Ontology Language, OWL. [5]. For those not familiar with description logics, we provide a very brief summary here covering the constructs used in SNOMED. (This is a subset of the fragment known as "EL++" or the "OWL EL profile.".) For more details, see The Description Logic Handbook [6] or the documentation and tutorials on the OWL web site. For readability, we shall use primarily the Manchester OWL Syntax [1], but the semantics are independent of the particular syntax. For convenience for readers who may use other references, we also give other equivalent and informal syntaxes.

Description logics are subsets of standard first order logic, and all statements can be translated into first order logic.

SNOMED "concepts" correspond to OWL "classes" and SNOMED attributes to OWL "properties". The hierarchical relation in SNOMED, otherwise known as sub-class/superclass relation, corresponds to "subsumption" in description logics and OWL.

In description logic, one class is subsumed by another if all of its members *necessarily* are members of the other, *i.e.* if it can be proven from the definitions and descriptions that any member of the subclass necessarily is a member of the superclass.

In fact, in description logics, the subclass-superclass relation is equivalent to necessary implication; to say that B is a subclass of A is to say that B necessarily implies A. Hence description logics have no separate operator for implication, such as "IF…THEN…" or "➔". They simply use "SubClassOf".

Although usually depicted as network, formally, a description logic model, or "ontology", consists of a set of axioms that always begin with "for all…". As used in SNOMED, there are three basic types.

1. **Simple subclass axioms between named classes: ("partial definitions")**
   examples
   *"Diabetes is a kind of Disorder"; "Pneumonia is a kind of Lung disease"*
   informally
   *"B is a kind of A"*
   which means
   *"All Bs are As",*
   or in OWL Manchester Syntax

*B SubclassOf A*[14]
       or in standard DL notation
$B \sqsubseteq A$
       or semi-formally
*for all x . IF B(x) THEN A(x)*
       in standard predicate logic notation
$\forall x. B(x) \rightarrow A(x)$

2. **Subclass axioms between a named subclass and a "restriction":**
       examples
*"Fingers are parts of hands", "Headache is located in the head",*
       informally
*Cs p Vs,* or in traditional Object-Attribute-Value notation *"C-p-V"*
       which means
*All Cs are linked by property p to some V*
       or in OWL Manchester syntax
*C SubClassOf ( p some V )*
       or in standard DL notation
$C \sqsubseteq \exists p . V$
       or semi-formally
*for all x . IF C(x) THEN there exists v . V(v) and p(x, v)*
       in standard predicate logic notation
$\forall x. C(x) \rightarrow \exists v . V(v)$ & $p(x,v).$

Note that in OWL, expressions such as "p some V" are known are "Existential restrictions". Strictly speaking they are classes – specifically the class of all individuals that are linked to some V by p. Hence the use of the SubClassOf keyword in OWL and the "$\sqsubseteq$" constructor in DL notation.

3. **Equivalence axioms between a named class and an expression ("Complete definitions")**
       examples
*"Heart disease" is defined as "Disorder that is located in some Heart structure"*;
*"Pneumococcal pneumonia" is defined as "Pneumonia that has causal agent some Pneumococcus"*
       informally
*Cs are logically equivalent to Expression_E,*
     which means
*Anything that satisfies Expression_E is a C, and all Cs satisfy Expression_E.*
       or in OWL Manchester syntax
*C EquivalentClass Expression_E*[15]
       or in standard DL notation
*C ≡ Expression_E*
       or semi-formally
*for all x . C(x) IF AND ONLY IF Expression_E(x)*
       in standard predicate logic notation
$\forall x . C(x) \leftrightarrow Expression\_E(x)$

In most description logics, including that used in SNOMED, the right-hand side of an axiom may be an arbitrary Boolean combination of class names and expressions, but the left-hand side must be a simple class name. In the description logic EL++ used in SNOMED, the Boolean expression on the right hand side is restricted to a series of nested conjunctions.

Classes that appear on the left hand side only of axioms of types 1 and 2 are called "partially defined" in SNOMED. We can say what is necessarily true of partially defined classes, but not how to recognise them. The classifier can infer that they should be classified under other concepts, but not that other concepts can be classified under them.

By contrast, classes that appear in at least one axiom of type 3 are known as "fully defined". Because they appear in an axiom that says that "Anything that satisfies the expression is a member of the class", the classifier

---

[14] sometimes abbreviated to "$B \rightarrow A$"
[15] often abbreviated $A \leftrightarrow Expression\_E$ or $A == Expression\_E$

can infer both that the class should be classified under other concepts and that other concepts should be classified under it.

Equivalence axioms for fully defined concepts act as rules for recognising other concepts that should be classified under them. This is how description logic classifiers are able to infer the classified hierarchy.

For example if we have the definitions and the axioms below:

> *"Heart disease" EquivalentClass ("Disorder" that site some "Heart structure")*
>
> *"Aortic stenosis" SubClassOf ("Disorder" that site some "Aortic valve")*
>
> *"Aortic valve" SubClassOf  "Heart structure"*

Then a classifier can prove that

> *"Aortic stenosis" SubClassOf "Heart disease"*

There are many different members of the class of description logics distinguished from each other by the operators they support.  In general, the more operators, the more difficult the proofs and the poorer the performance of automatic classifiers.  The description logic used by SNOMED is known as "EL++"; the corresponding profile of OWL is known as "OWL-EL". It has been chosen so that classification is highly efficient even for very large models. It contains only the operators in the examples above.  It does not support negation or disjunction or expressions containing "only" (universal restrictions).  Algorithms have recently been found that are particularly efficient for classifying models formulated in EL++ [7].  SNOROCKET is based on these algorithms.  The field of classifier development for more expressive supersets of EL++ is developing rapidly, but none of the experimental systems is yet fully deployed at the time of writing.

In summary, key points about description logic semantics are:

1. For one concept to be a child (or descendant) of another means that all instances of the descendant are necessarily also instances of the ancestor.

2. All statements are universal, *i.e.* begin with "for all…".  Therefore, anything said about a concept is said about all of its instances.

3. Combining 1) and 2), all statements made about a concept in a description logic apply to all its children and descendants.

## Appdendix VI.      Tools used and availability

### VI.1.   SNOMED stated form

The SNOMED stated form is part of the standard distribution of SNOMED under …/OtherResources/StatedRelationships… A Perl script for converting the given form is provided that contains instructions for its use.  The availability and mechanism for obtaining the SNOMED distribution depends on the arrangements for each country, but is available worldwide for academic use.  Contact http://www.ihtsdo.org.

### VI.2.   Editors and visualisation tools

Protégé 4.0 and 4.1 available from http://protégé.stanford.org.  At the time of writing (Dec 2010) the conversion between version 4.1 and 4.2 necessitated by late changes to the OWL-2 specification had not been completed for all plugins.  Therefore, an awkward combination of the two versions was required.  In particularly, SNOROCKET was only available with PROTÉGÉ 4.0 and OPPL with Protégé 4.1.  A special modification of Protégé 4.0 was used that keeps class names in alphabetical order, which greatly eases analysis. It is available from the Protégé website.

Protégé supports many plugins and views.  Essential to this work are the OWLViz graphical view, the inferred superclass view, the usage view and the query view. (See protégé documentation for information about installing views)

Note the Perl script produces an OWL file with class names of the form SCT_snomed_id and the SNOMED fully specified name as the label annotation.  In order to view the file using the fully specified names there, the preferences must be set so that "Renderer" is by "Annotation".  Also note that because the SNOMED names contain spaces, most searches, all expressions in editors and all queries must begin with a single quote. Since SNOMED fully specified names are used, the auto-complete facility (tab or control-space) is particularly helpful in adding the suffixes and ending single quotes.

The SNOMED ID can be obtained by pressing control-U (CMD-U on Mac) to bring up the actual name

Protégé 4.1 supports "justifications" – explanations for unexpected inferences.  If the root cause of an inference is unclear, justifications can usually help to identify it.  The justifications are activated by clicking the circled question mark beside any inferred (shaded) class or restriction.

### VI.3.    Browsers for full release of SNOMED

Besides Protégé 4 browser, we used primarily SNOB developed by Egbert Van Der Hering following work on the GALEN project (http://snob.eggbird.eu/)   SNOB provides an easier way to browse the inferred superclass hierarchy than other browsers.  However, the standard CliniClue Xplore (http://www.cliniclue.com)  and IHTSDO stand alone browser (available from IHTSDO – http://www.ihtsdo.org) were also used to validate results.

### VI.4.    Module extraction

The modularization tools are built into the new OWL API 3 (confusingly the complete API for OWL-2) available from http://www.owlapi.sourceforge.net.  A simple Java program giving direct access to the feature along with instructions in its use can be found at http://owl.cs.manchester.ac.uk/snomed/.  It also includes instructions for extracting subsets as "signatures" for modules using Protégé.  It is expected that these tools will be updated, but the updates will remain available from this site.

### VI.5.    Classifiers

The main classifier used was SNOROCKET available from http://aehrc.com/hie/snorocket.html.  (NB, although SNOROCKET is available as a plugin for Protégé it is not available through the standard Protégé plugin library but must be downloaded separately.) The three other classifiers available with Protégé were also used for validation and cross checking, FaCT++, Pellet, and Hermit.  All three are available as plugins under the preferences section of Protégé. (Pellet requires a separate license for commercial use, although different versions will be downloaded with Protégé 4.0 and 4.1.)

### VI.6.    Combined lexical and semantic searching plus scripting & testing

For sophisticated search (and repair) the Ontology Pre-Processing Language (OPPL), which is available as a plugin for Protégé 4.1.  OPPL can be used either to formulate complex queries or to execute scripts based on those queries.  It can also be integrated with the Protégé Unit Testing Framework.  (See http://oppl2.sourceforge.net/)

## Appdendix VII.    Additional screenshots documenting unexpected classifications

The standard tools for viewing SNOMED often obscure the upwards hierarchies.  Below are supplementary screenshots from the SNOB browser on the Jan 31 2010 release of SNOMED documenting assertions in the main text.

### VII.1.    Autoimmune and Allergic disorders

Standard hierarchy showing three misplaced autoimmune diseases:

- Grave's disease
- Myasthenia gravis
- Idiopathic thrombocytopenic purpura

### VII.2.   Myocardial infarction – not classified under Ischemic heart disease

```
[+] Injury of anatomical site
  [↓] Heart disease
    [+] Myocardial finding
[−] Myocardial disease
  [+] Heart disease
  [−] Structural disorder of heart
[ ] Myocardial infarction
```
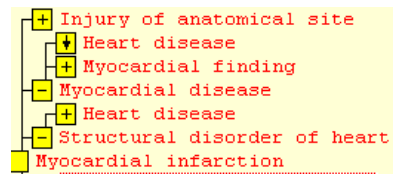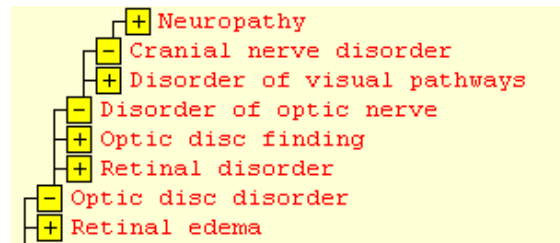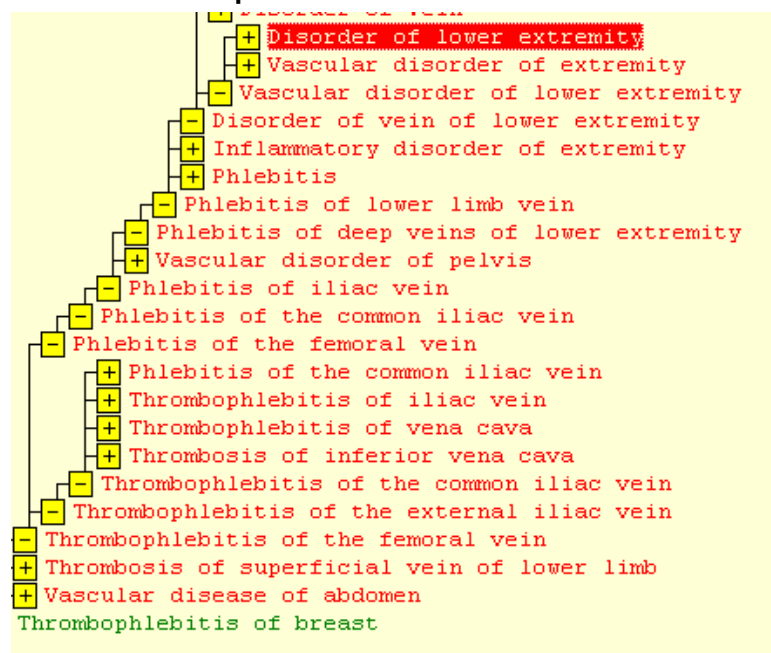
### VII.3.   Optic disc edema and neuropathy

```
        [+] Neuropathy
      [−] Cranial nerve disorder
        [+] Disorder of visual pathways
    [−] Disorder of optic nerve
      [+] Optic disc finding
      [+] Retinal disorder
  [−] Optic disc disorder
    [+] Retinal edema
```

### VII.4.   Thrombophlebitis of the breast and relation to abdomen and lower limb

```
          [+] Disorder of lower extremity
          [+] Vascular disorder of extremity
        [−] Vascular disorder of lower extremity
       [−] Disorder of vein of lower extremity
        [+] Inflammatory disorder of extremity
        [+] Phlebitis
      [−] Phlebitis of lower limb vein
     [−] Phlebitis of deep veins of lower extremity
       [+] Vascular disorder of pelvis
   [−] Phlebitis of iliac vein
  [−] Phlebitis of the common iliac vein
 [−] Phlebitis of the femoral vein
     [+] Phlebitis of the common iliac vein
     [+] Thrombophlebitis of iliac vein
     [+] Thrombophlebitis of vena cava
     [+] Thrombosis of inferior vena cava
   [−] Thrombophlebitis of the common iliac vein
  [−] Thrombophlebitis of the external iliac vein
[−] Thrombophlebitis of the femoral vein
[+] Thrombosis of superficial vein of lower limb
[+] Vascular disease of abdomen
Thrombophlebitis of breast
```

### VII.5.   Subdural hemorrhage

The full OWLViz visualization of subdural hemorrhage is too large to be held in a MSWord document and is provided as a .jpg image.

## Appendix References:

1.    Horridge M, Patel-Schneider PF. OIWL 2 Web Ontology language: Manchester Syntax. 2009; http://www.w3.org/TR/owl2-manchester-syntax/.

2.    Rector AL, Bechhofer S, Goble CA, Horrocks I, Nowlan WA, Solomon WD. The GRAIL concept modelling language for medical terminology. Artificial Intelligence in Medicine. 1997;9:139-171.

3.    Horrocks I; Using an expressive description logic: FaCT or Fiction. 1998; Principles of Knowledge Representation and Reasoning: Proceedings of the Sixth International Conference on Knowledge Representation (KR 98): San Francisco, CA: Morgan Kaufmann; 634-647.

4.    Horrocks I, Sattler. U. The decidability of SHIQ with complex role inclusion axioms. Artificial Intelligence. 2004;160:79-104.

5.    W3C OWL Working Group. OWL2 Web Ontology Language Document Overview. http://www.w3.org/TR/owl2-overview/.

6.    Baader F, Calvanese D, McGinness DL, Nardi D, Patel-Schneider PF, editors. The Description Logic Handbook. Cambridge, England: Cambridge University Press; 2003:555.

7.    Baader F, Brandt S, Lutz C; Pushing the EL Envelope. 2005; Proceedings of the Nineteenth International Joint Conference on Artificial Intelligence IJCAI-05: 364–369.